

# Self-Supervised Representation Learning



Raul C. Incze

- Motivation
  - Sample Efficiency
- Representations
  - The human brain
  - The machine
  - The properties of a good representation
- Self-Supervision
  - Autoencoding (AE) Methods
  - Autoregressive (AR) Methods
  - Pretext Tasks
- Applications
  - Few-Labels Scenarios
  - Reinforcement Learning

# A little about myself...

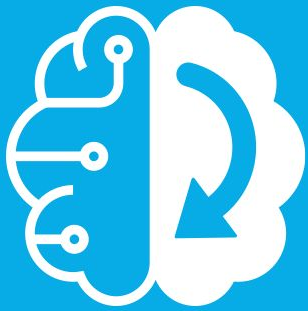
## Education

- BSc, Computer and Information Sciences - Babeş-Bolyai University (2015)
- MSc, Applied Computational Intelligence - Babeş-Bolyai University (2017)

## Work Experience

- Game Developer (4 years)
- Machine Learning Engineer (3 years)





**cognifeed**

# Cogni-what?

Accessible Machine Teaching:

- Empowers people with the domain knowledge to train models.
- Sample and label efficiency.

# Cogni-what?

Accessible Machine Teaching:

- Empowers people with the domain knowledge to train models.
- **Sample and label efficiency.**
  - Creating new datasets for each task is expensive.
  - Some domains are supervision-starved.
  - There are many unlabelled data samples.

$$\frac{1}{\varepsilon(1 - \sqrt{\varepsilon})} [2d \ln(6/\varepsilon) + \ln(2/\delta)]$$

- VC dimension ~ the effective number of parameters (expressiveness).
- In practice, the number of samples = 10 x VC dimension.<sup>2</sup>
- Deep nets => huge VC dimension

1. [Bounding sample size with the VC dimension](#), Shawe-Taylor, J. et al, 1993  
2. [The VC Dimension - A measure of what it takes a model to learn](#), Yaser Abu-Mostafa, 2012

# Conditioning is important

- Priors
- Regularization

$$\tilde{O} \left( (m + r) / \varepsilon^2 \right)$$



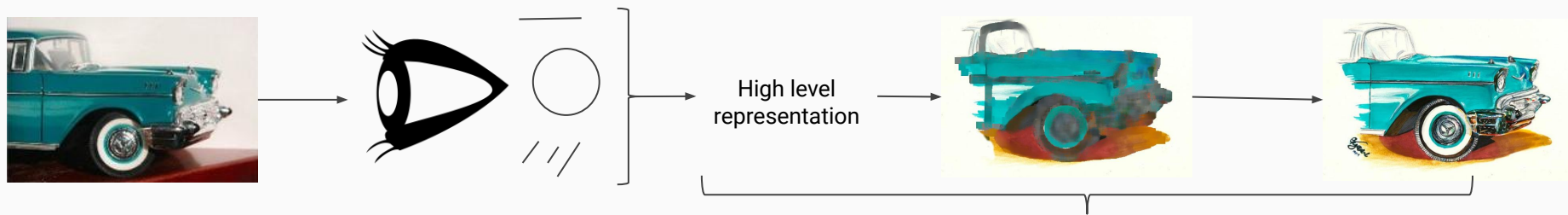
1. [The Model Complexity Myth](#), Jake VanderPlas, 2015
2. [How Many Samples are Needed to Estimate a Convolutional Neural Network](#), Simon S. Du et al, 2018



We can do better...

# Representations

# The Human Brain and Representation Learning



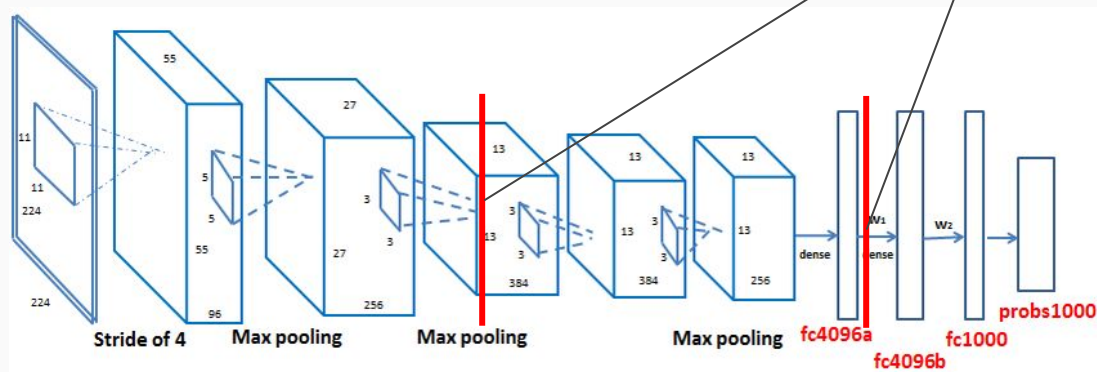
Some takeaways:

- We store representations, not snapshots.
- We use these representations to reason (predict, classify) new inputs and recall old ones.
- We don't need "supervision" to generate representations.



# Representations in Machine Learning

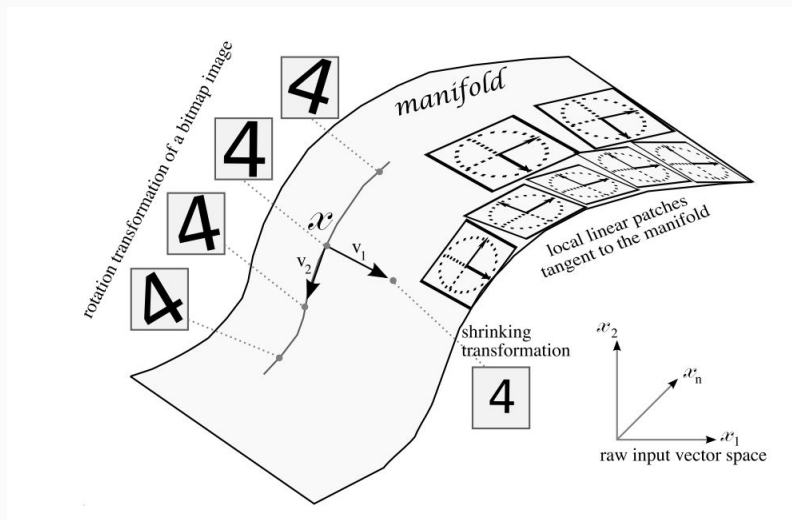
- Also known as “feature learning”.
- It’s a **consequence** of deep supervised learning.



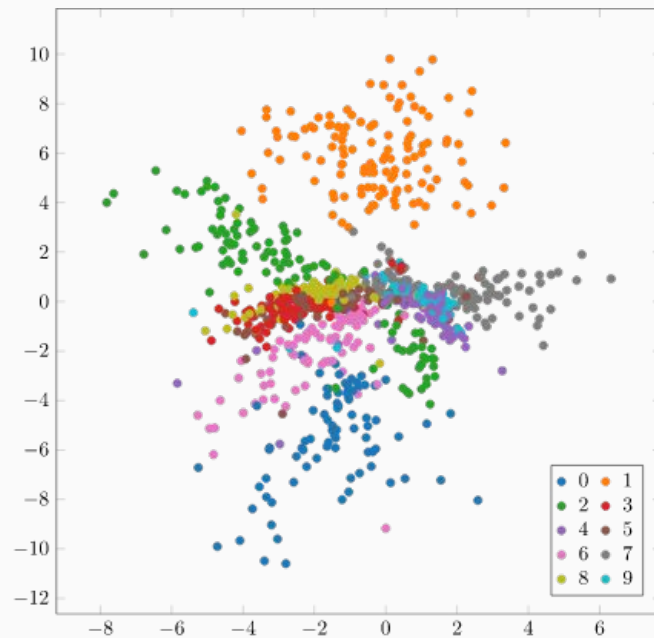
We can also make models learn representations **explicitly**.

# What makes a good representation?

- **Smoothness**
- Multiple explanatory factors
- Depth
- **Shared factors across tasks**
- **Manifolds**
- **Natural clustering**
- **Temporal and spatial coherence**
- Sparsity
- Simplicity of factor dependencies



Smoothness



Natural clustering

# Self-Supervision

# Self-supervision?

- A form of unsupervised learning.
- Data itself provides the supervision.

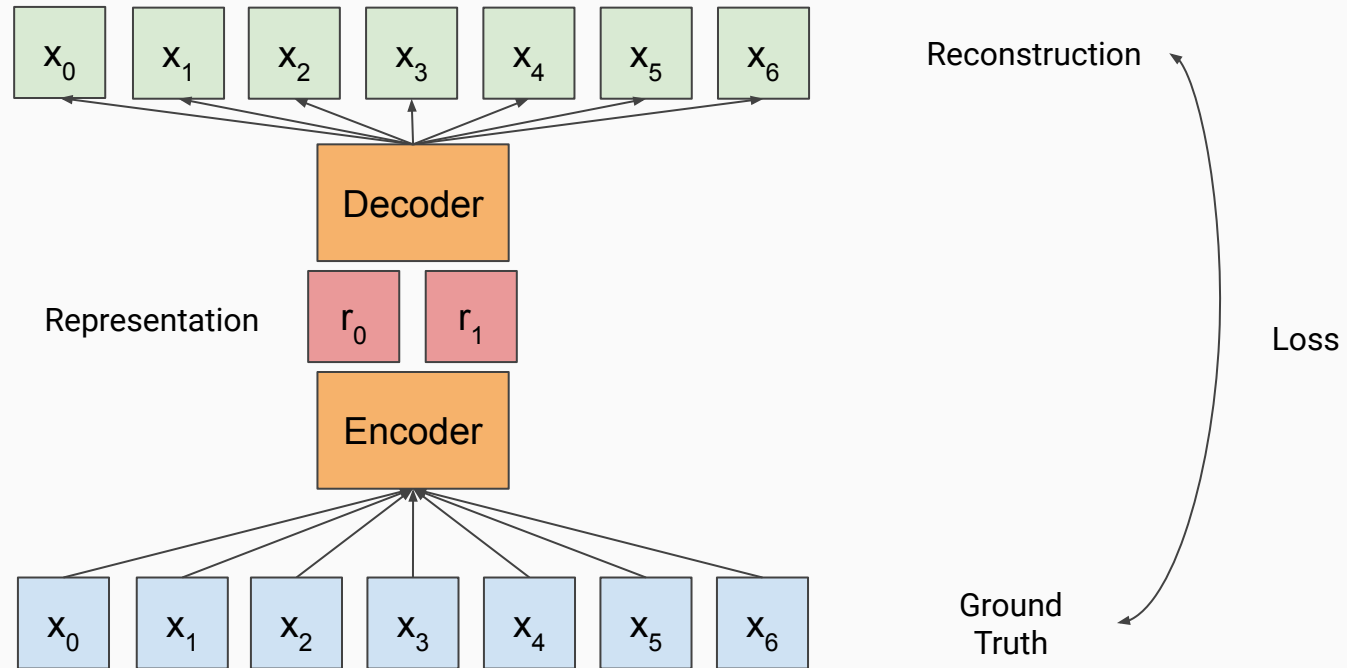
Two ways of achieving this:

- We withhold part of the input data and make our model predict it.
- Making up pretext tasks.



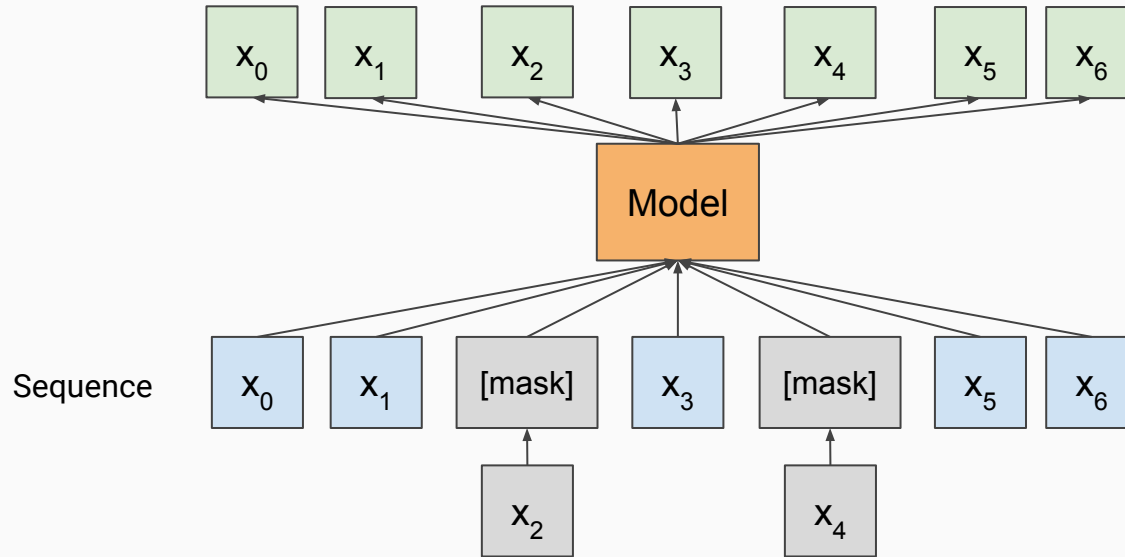
# Withholding Data - Autoencoding (AE) Methods

# Autoencoding (AE) Methods



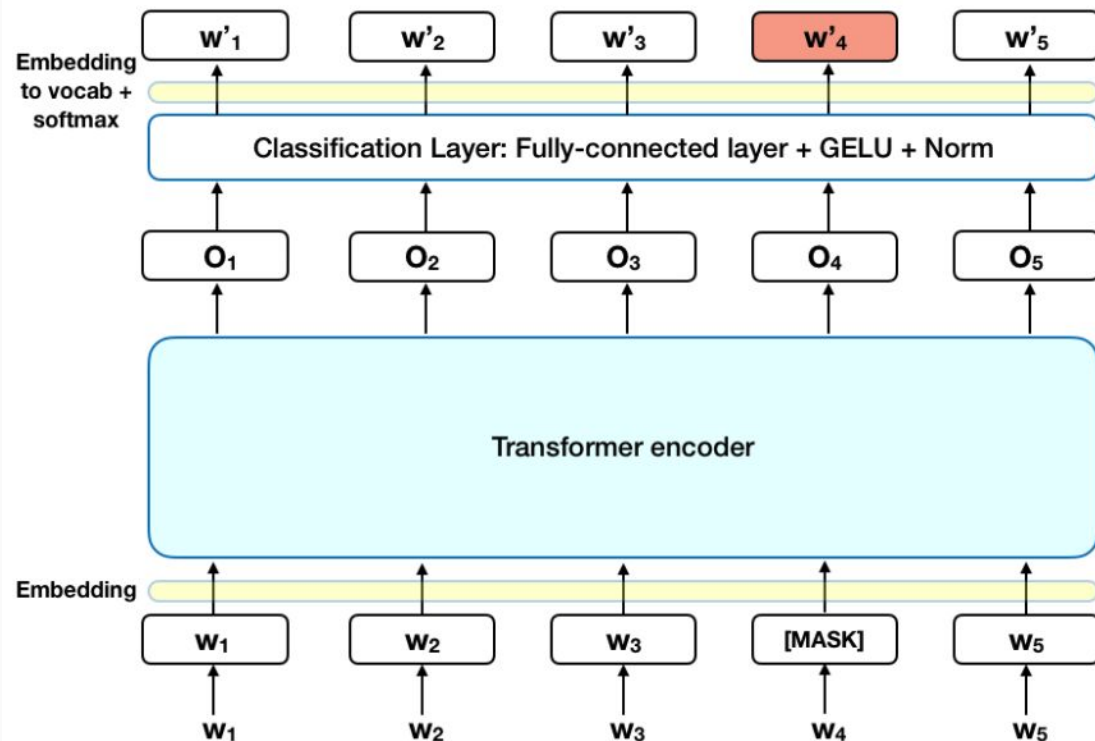
No information withheld (or arguably all the input is withheld from the decoder), representation is quite arbitrary, reconstruction can be poor.

# Autoencoding (AE) Methods



- Denoising autoencoder (noise = input values set to 0/mask value).
- Works under the assumption that the choice for  $x_2$  is independent from the choice of  $x_4$ .
- Natural Language breaks this assumption but this still works well on NLP tasks.
- Model = Transformer => BERT

# Study Case: Bidirectional Encoder Representations from Transformers (BERT)



## Fine-tuning approach

$BERT_{LARGE}$	96.6
$BERT_{BASE}$	96.4

## Feature-based approach ( $BERT_{BASE}$ )

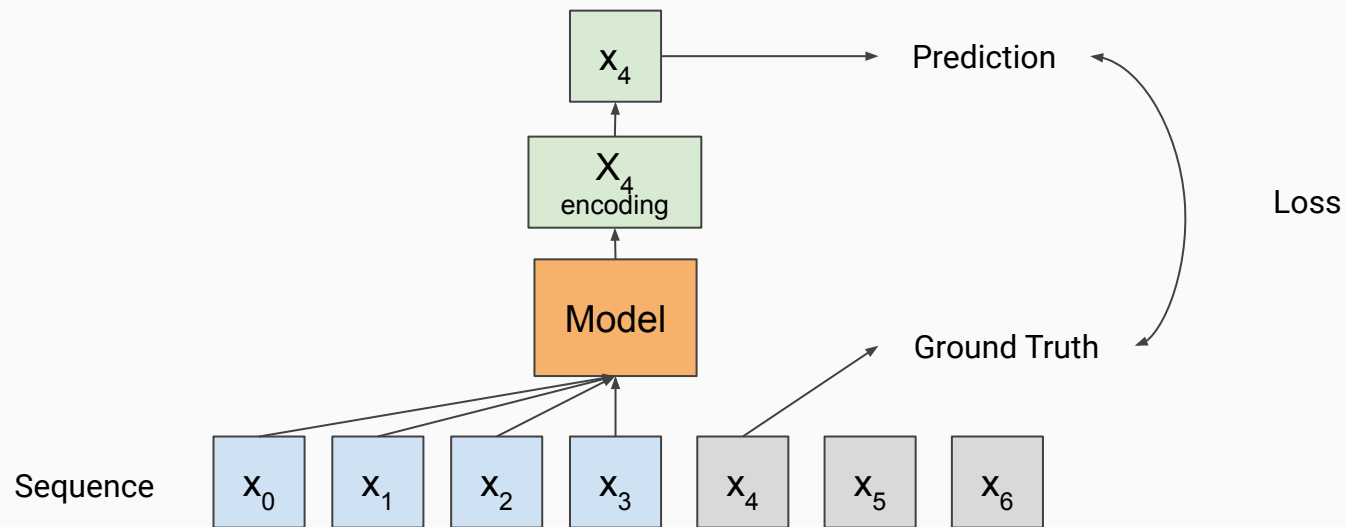
Embeddings	91.0
Second-to-Last Hidden	95.6
Last Hidden	94.9
Weighted Sum Last Four Hidden	95.9
Concat Last Four Hidden	96.1
Weighted Sum All 12 Layers	95.5

Results on Named Entity Recognition

1. [Attention Is All You Need](#), Ashish Vaswani et. al, 2017
2. [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#), Jacob Devlin et. al, 2018

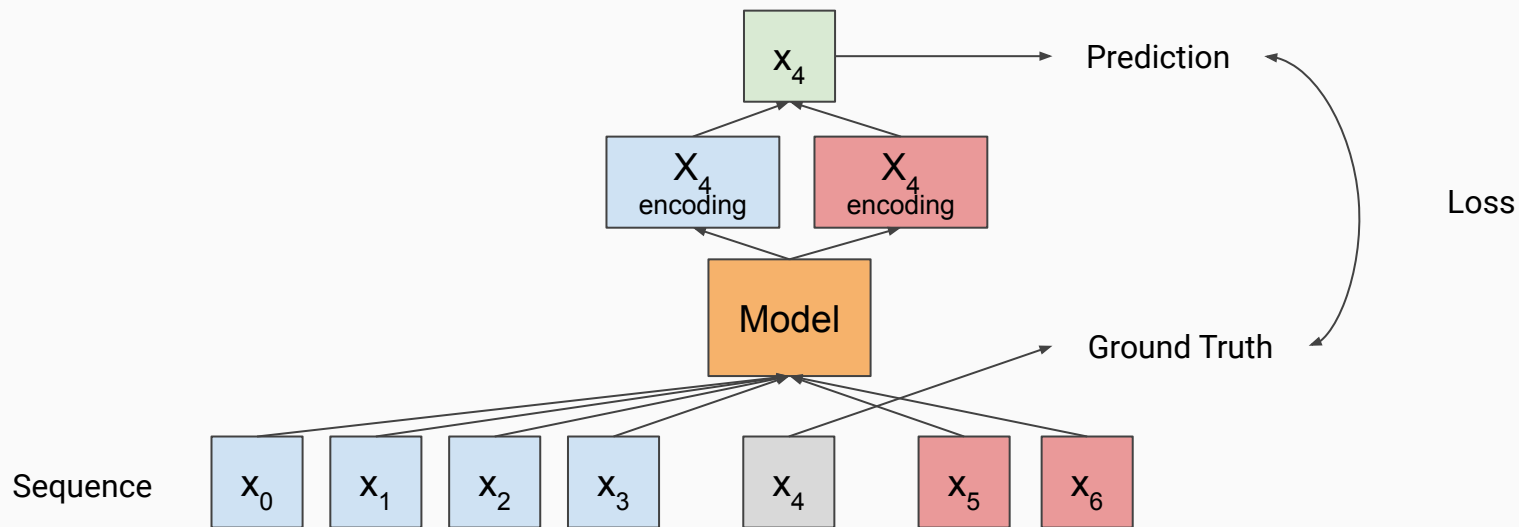
# Withholding Data - Autoregressive (AR) Methods

# Autoregressive Methods



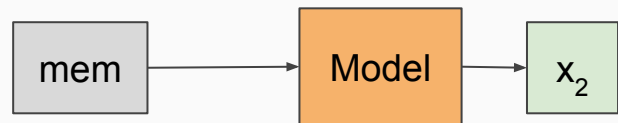
- Assumes that for the prediction of  $x_t$  all we need is  $x_0, \dots, x_{t-1}$
- Model = Transformer  $\Rightarrow$  GPT (OpenAI language model).

# Autoregressive Methods



- Assumes the whole bidirectional context.
- Model = LSTM => ELMO.

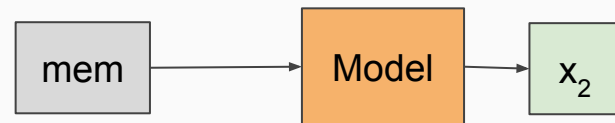
# Autoregressive Methods



Sequence



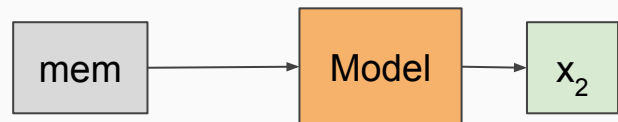
Factorization order:  $2 \square 1 \square 3 \square 0$



Sequence



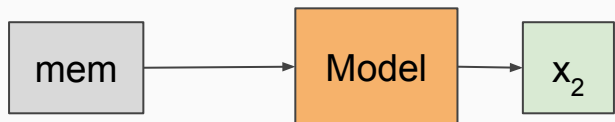
Factorization order:  $1 \square 3 \square 2 \square 0$



Sequence



Factorization order:  $0 \square 3 \square 1 \square 2$



Sequence

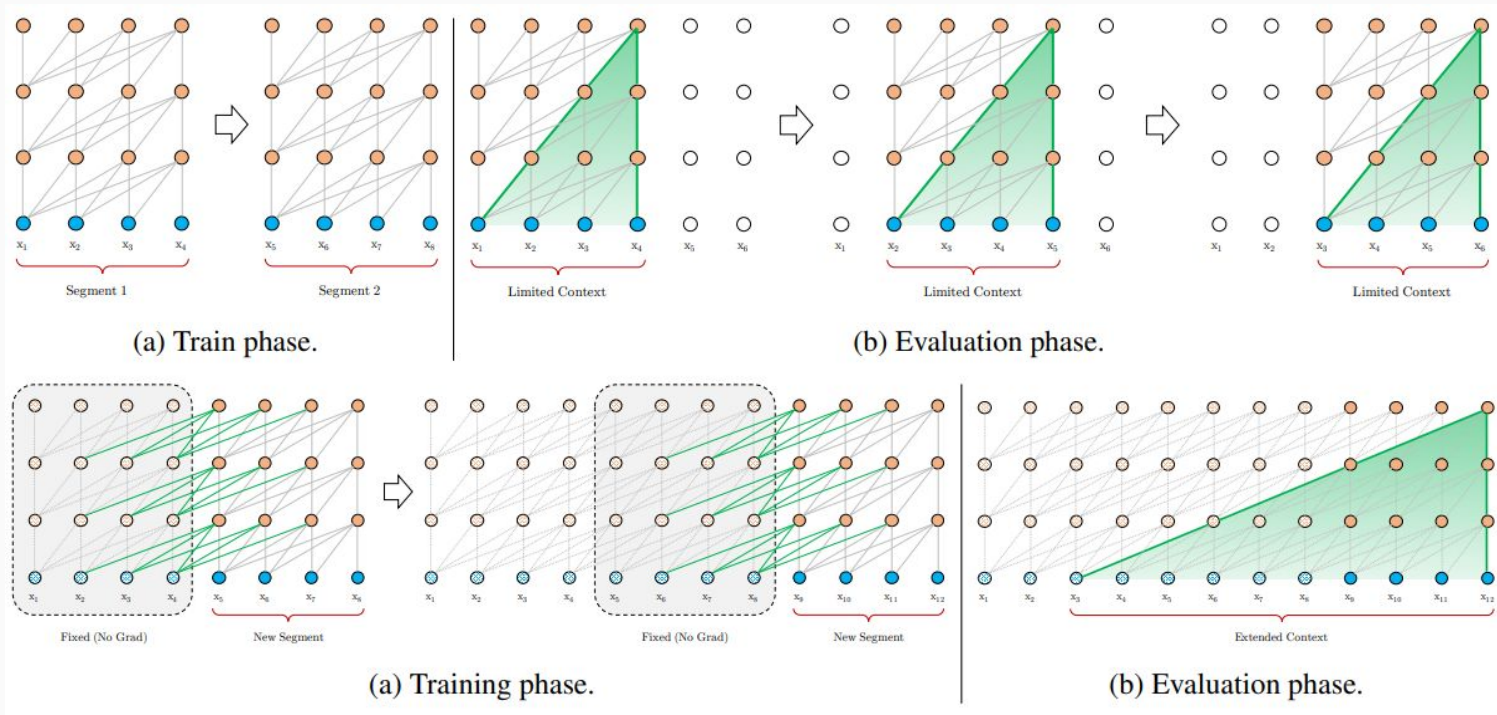


Factorization order:  $3 \square 2 \square 0 \square 1$



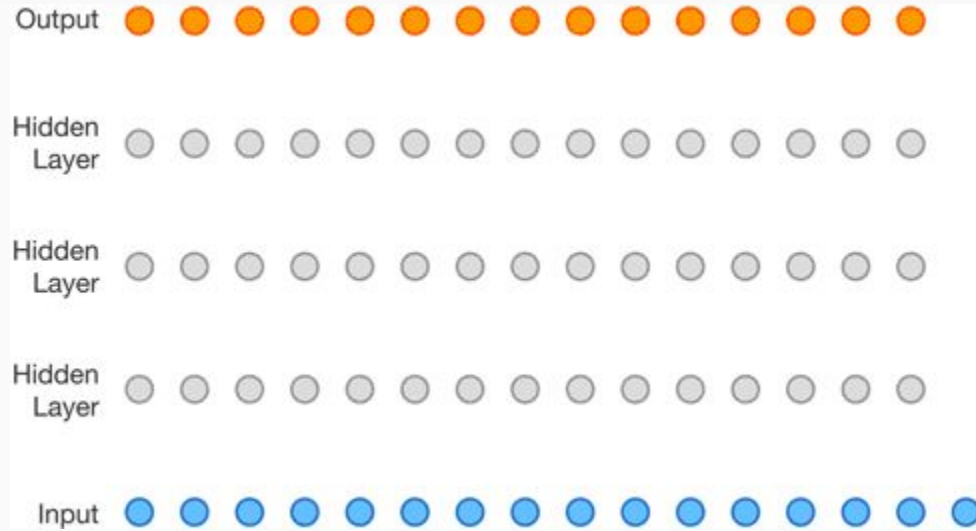
# Case Study: XLNet

## Factorization + Transformer-XL

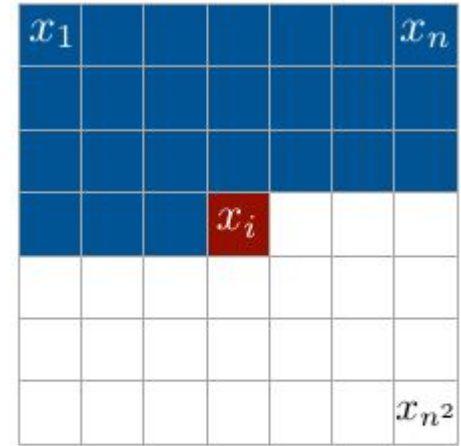


1. [Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context](#), Zihang Dai et. al, 2019
2. [XLNet: Generalized Autoregressive Pretraining for Language Understanding](#), Zhilin Yang et. al, 2019

# Other Autoregressive Models



WaveNet



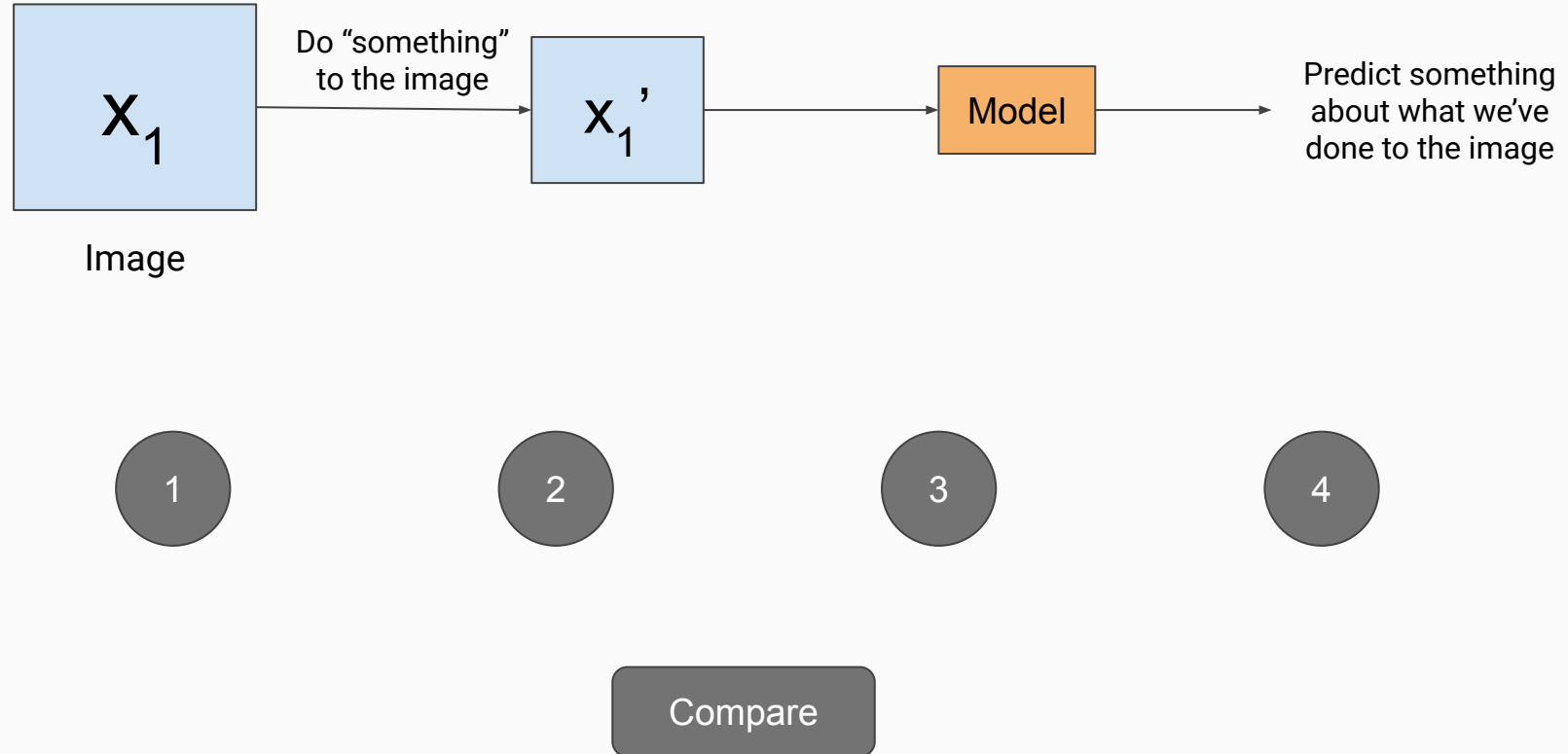
Context

PixelRNN

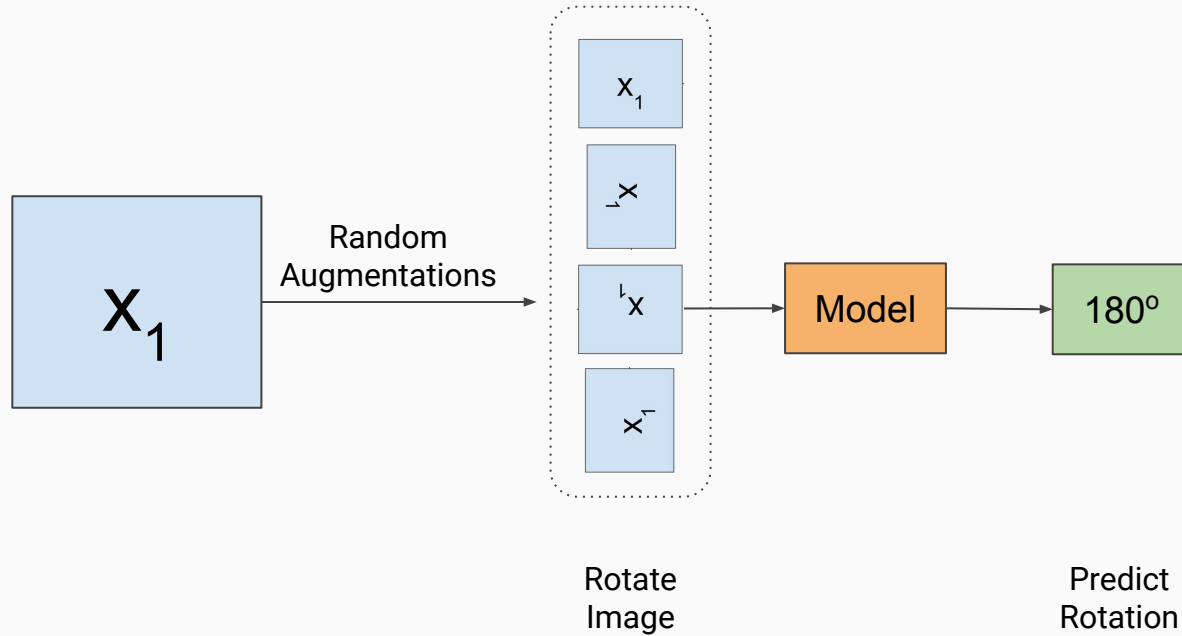
1. [WaveNet: A Generative Model for Raw Audio](#), Aaron van den Oord et. al, 2016
2. [Pixel Recurrent Neural Networks](#), Aaron van den Oord et. al, 201

# Pretext Tasks

# Brainstorming Exercise - Pretext Tasks For Visual Representation Learning

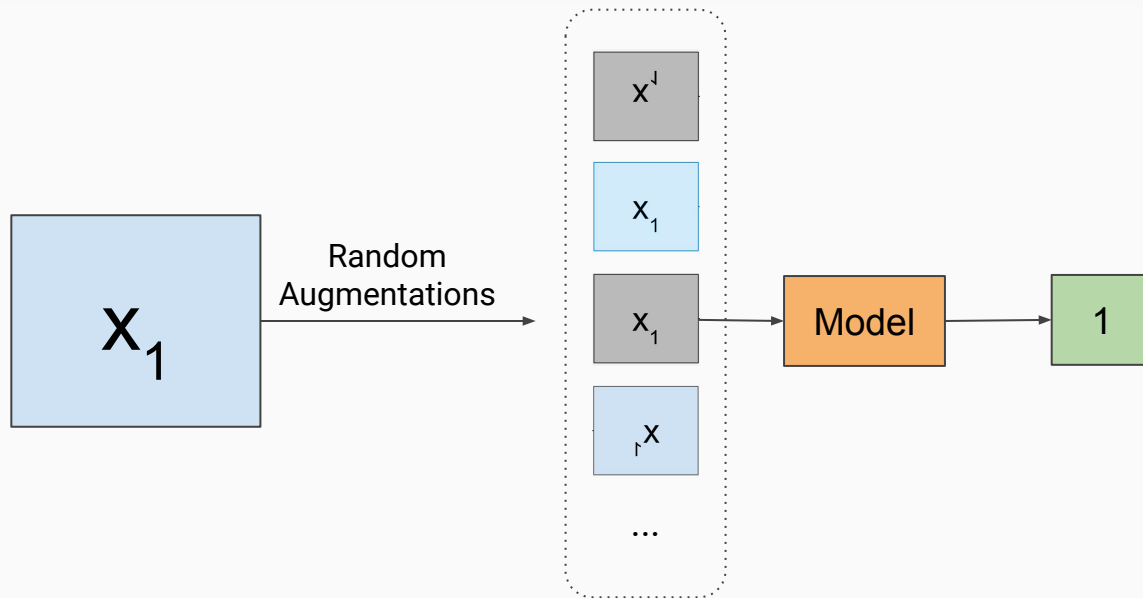


# Rotation



Back

# Exemplar

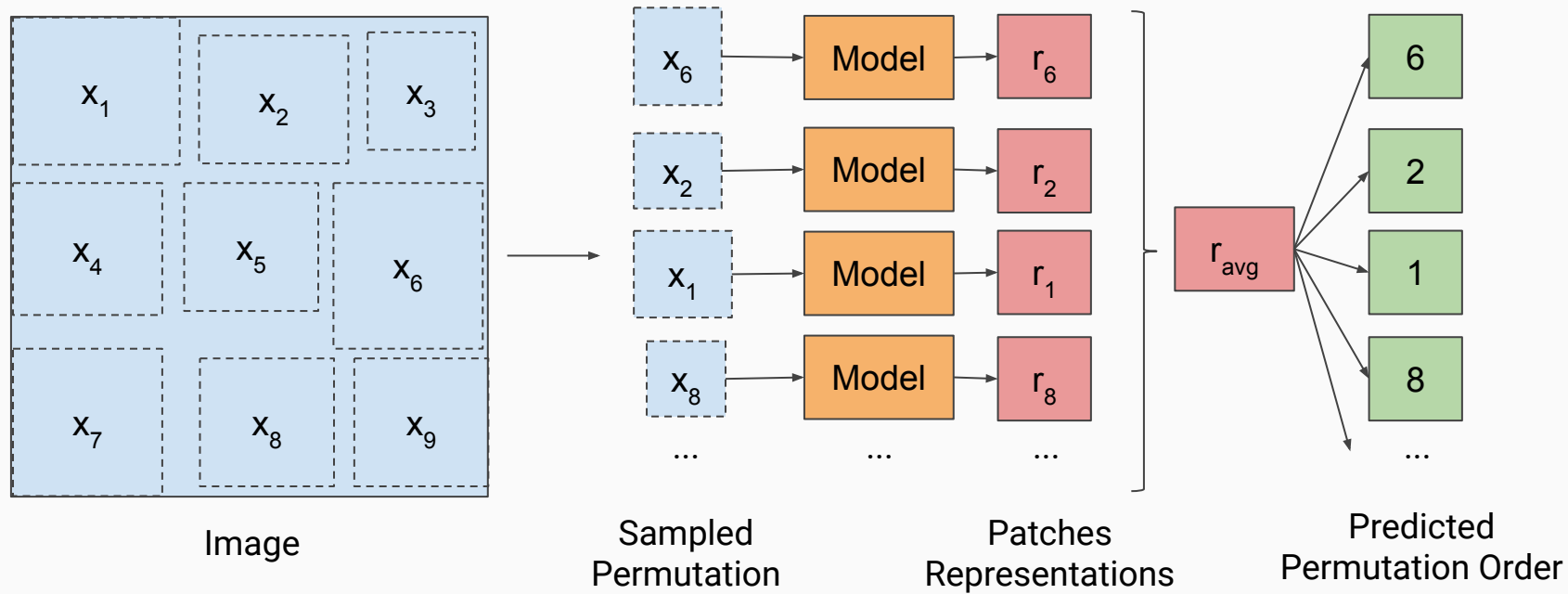


Every original image is an "exemplar" and has a corresponding class

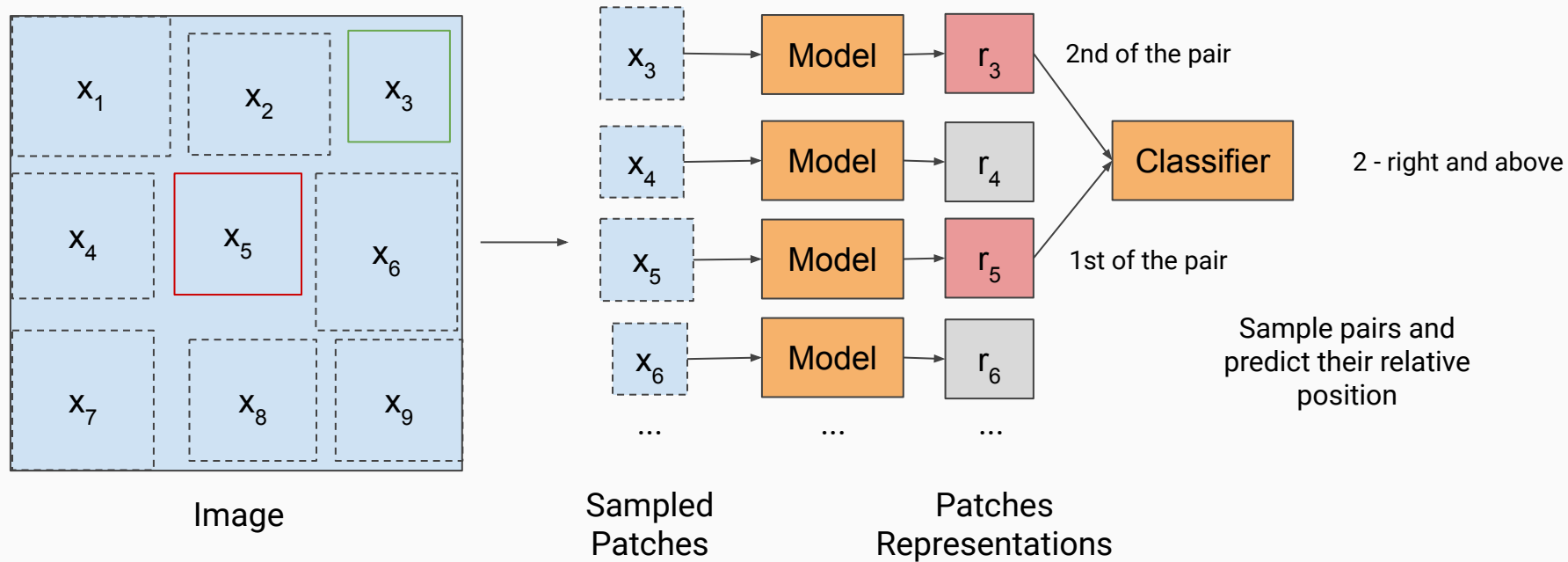
Label all augmentations of  $x_1$  with  $x_1$ 's exemplar label "1"

Model has to predict  $x_1$  exemplar label

Back



# Relative Patch Location





# Comparison

Model	Rotation				Exemplar			RelPatchLoc		Jigsaw	
	4×	8×	12×	16×	4×	8×	12×	4×	8×	4×	8×
RevNet50	<b>47.3</b>	<b>50.4</b>	<b>53.1</b>	<u><b>53.7</b></u>	<b>42.4</b>	45.6	46.4	40.6	45.0	40.1	43.7
ResNet50 v2	43.8	47.5	47.2	<u>47.6</u>	<b>43.0</b>	45.7	46.6	42.2	46.7	38.4	41.3
ResNet50 v1	41.7	43.4	43.3	43.2	<b>42.8</b>	<b>46.9</b>	<b>47.7</b>	<b>46.8</b>	<u><b>50.5</b></u>	<b>42.2</b>	<b>45.4</b>
RevNet50 (-)	45.2	<b>51.0</b>	<b>52.8</b>	<u><b>53.7</b></u>	38.0	42.6	44.3	33.8	43.5	36.1	41.5
ResNet50 v2 (-)	38.6	44.5	47.3	<u>48.2</u>	33.7	36.7	38.2	38.6	43.4	32.5	34.4
VGG19-BN	16.8	14.6	16.6	22.7	26.4	28.3	<u>29.0</u>	28.5	<u>29.4</u>	19.8	21.1

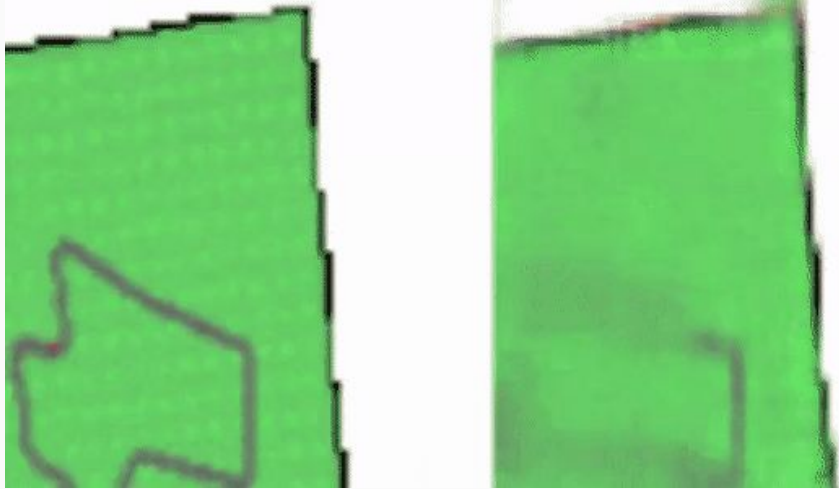
1. [Revisiting Self-Supervised Visual Representation Learning](#), Alexander Kolesnikov et. al, 2019

# Putting Representations to Good Use

# Representations In Few-Labels Scenarios

- Pretrain on vast amounts of data in a self-supervised manner and then:
  - Fine-tune on known labels.
  - Use learned representations to train a shallow model.
- Use representations' properties (natural clustering) for label propagation.
- Active Learning
  - Representations are learned.
  - Estimator is trained on available labels.
  - Estimator predicts labels for unlabeled instances.
  - Forwards the most uncertain instances to the user to be corrected.
- Co-Training
  - Learn two different representations for each data point => two views.
  - Train an estimator on each view.
  - Each estimator predicts labels for unlabeled instances and adds the most confident prediction to the label pool of the other estimator.

# Self-Supervised Representations in Reinforcement Learning



Agents learning in their own “dreams”<sup>1</sup>



(a) learn to explore on Level-1



(b) explore faster on Level-2

Curiosity-driven exploration<sup>2</sup>

1. [World Models](#), David Ha and JÜRGEN SCHMIDHUBER, 2018
2. [Curiosity-driven Exploration by Self-supervised Prediction](#), Deepak Pathak et. al, 2017

Question Time